

Chapter 2

Optimization

Gradients, convexity, and ALS



Contents

- Background
- Gradient descent
- Stochastic gradient descent
- Newton's method
- Alternating least squares
- KKT conditions

Motivation

- We can solve basic least-squares linear systems using SVD
- But what if we have
 - missing values in the data
 - extra constraints for feasible solutions
 - more complex optimization problems (e.g. regularizers)
- etc

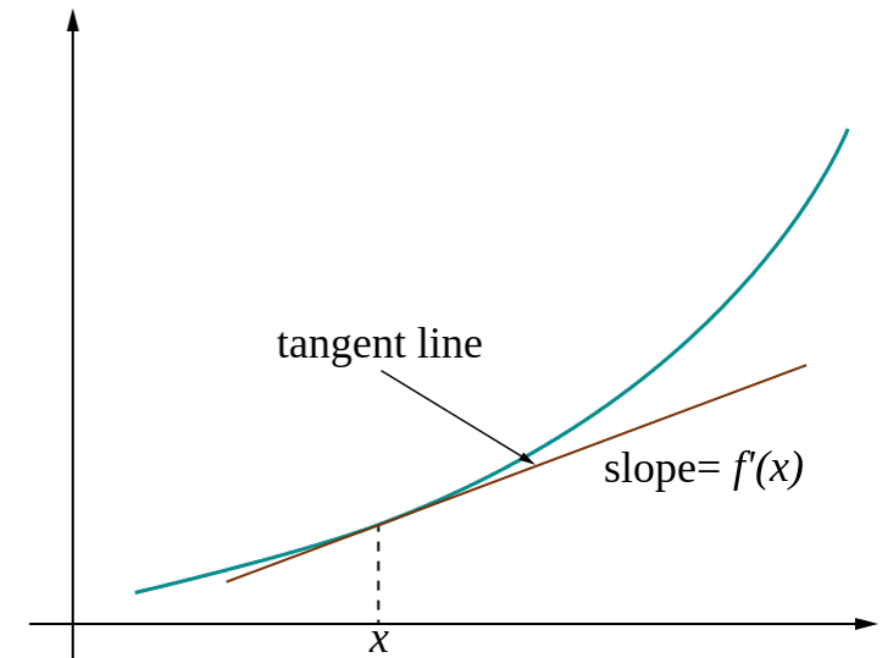
Gradients, Hessians, and convexity

Derivatives and local optima

- The **derivative** of a function $f: \mathbb{R} \rightarrow \mathbb{R}$, denoted f' , explains its *rate of change*

$$f'(a) = \lim_{h \rightarrow 0^+} \frac{f(a+h) - f(a)}{h}$$

- If it exists
- The second derivative f'' is the change of rate of change



Derivatives and local optima

- A **stationary point** of differentiable f is x s.t.
 $f'(x) = 0$
 - f achieves its extremes in stationary points or in points where derivative doesn't exist, or at infinities (Fermat's theorem)
- Whether this is (local) maximum or minimum can be seen from the second derivative (if it exists)

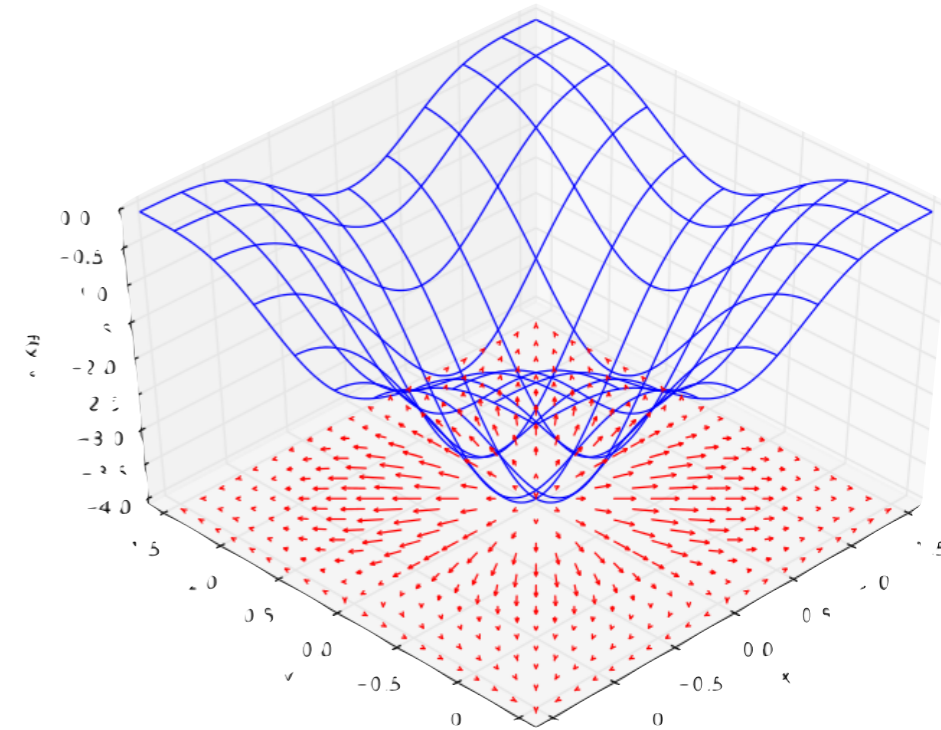
Partial derivative

- If f is multivariate (e.g. $f: \mathbb{R}^3 \rightarrow \mathbb{R}$), we can consider it as a family of functions
 - E.g. $f(x, y) = x^2 + y$ has functions $f_x(y) = x^2 + y$ and $f_y(x) = x^2 + y$
- **Partial derivative** w.r.t. one variable keeps other variables constant

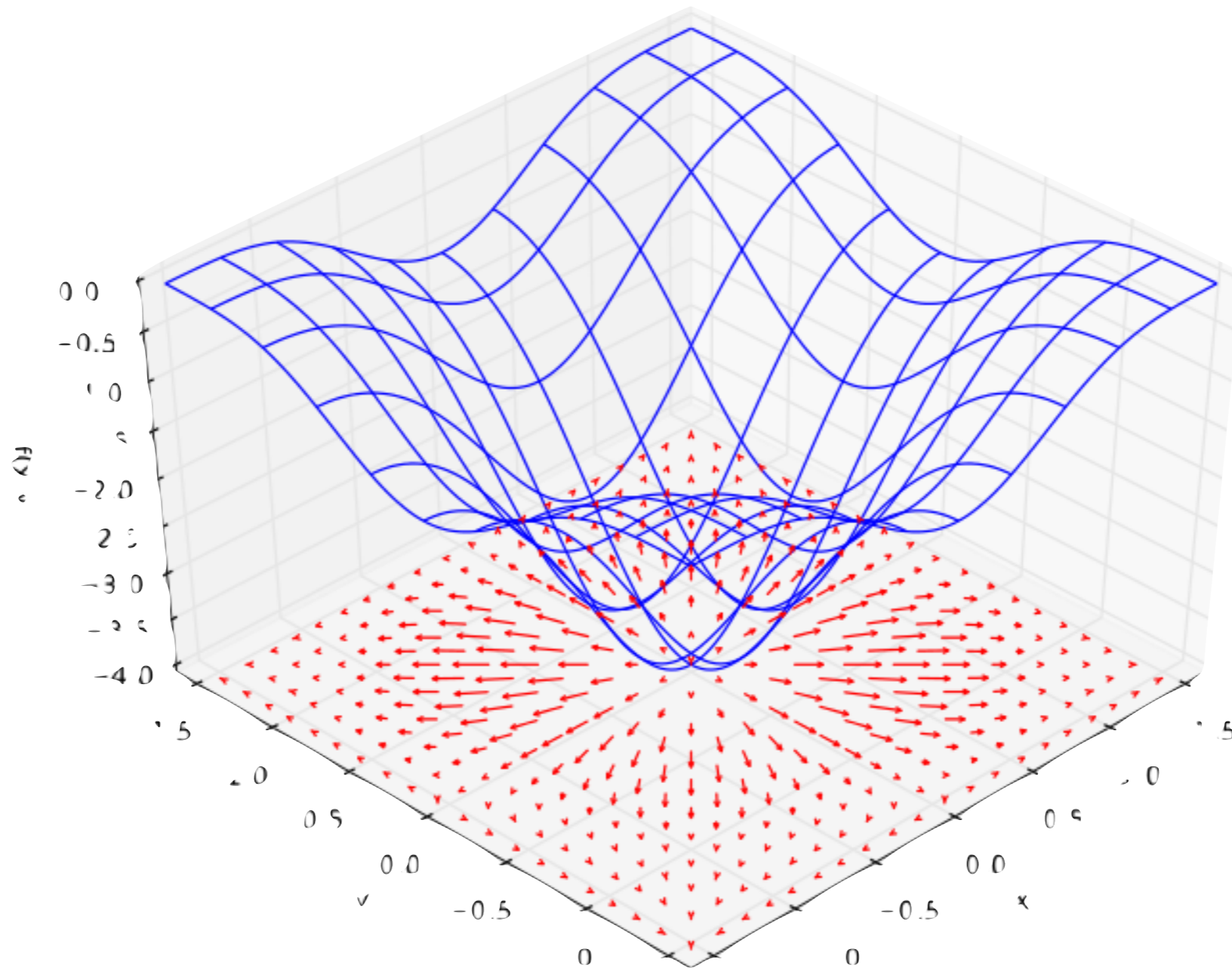
$$\frac{\partial f}{\partial x}(x, y) = f'_y(x) = 2x$$

Gradient

- **Gradient** is the derivative for multivariate functions $f: \mathbb{R}^n \rightarrow \mathbb{R}$
- $\nabla f = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$
- Here (and later), we assume that the derivatives exist
- Gradient is a function $\nabla f: \mathbb{R}^n \rightarrow \mathbb{R}^n$
- $\nabla f(\mathbf{x})$ points “up” in the function at point \mathbf{x}



Gradient



Hessian

- **Hessian** is a square matrix of all second-order partial derivatives of a function

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\mathbf{H}(f) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

- As usual, we assume the derivatives exist

Jacobian matrix

- If $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$, then its **Jacobian** (matrix) is an $n \times m$ matrix of partial derivatives in form

$$\mathbf{J} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_m} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \cdots & \frac{\partial f_n}{\partial x_m} \end{pmatrix}$$

- Jacobian is the best linear approximation of f
- $\mathbf{H}(f(\mathbf{x})) = \mathbf{J}(\nabla f(\mathbf{x}))^T$

Examples

Function

$$f(x, y) = x^2 + 2xy + y$$

Partial derivatives

$$\frac{\partial f}{\partial x}(x, y) = 2x + 2y$$

$$\frac{\partial f}{\partial y}(x, y) = 2x + 1$$

Gradient

$$\nabla f = (2x + 2y, 2x + 1)$$

Hessian

$$\mathbf{H}(f) = \begin{pmatrix} 2 & 2 \\ 2 & 0 \end{pmatrix}$$

Function

$$f(x, y) = \begin{pmatrix} x^2 y \\ 5x + \sin y \end{pmatrix}$$

Jacobian

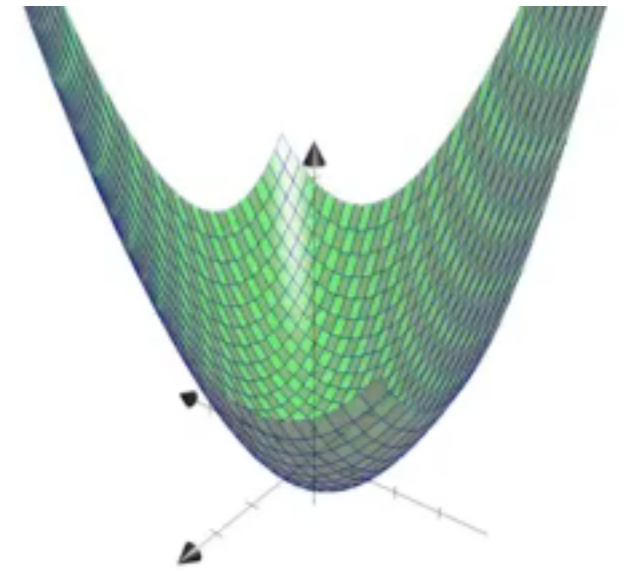
$$\mathbf{J}(f) = \begin{pmatrix} 2xy & x^2 \\ 5 & \cos y \end{pmatrix}$$

Gradient's properties

- **Linearity:** $\nabla(\alpha f + \beta g)(\mathbf{x}) = \alpha \nabla f(\mathbf{x}) + \beta \nabla g(\mathbf{x})$
- **Product rule:** $\nabla(fg)(\mathbf{x}) = f(\mathbf{x})\nabla g(\mathbf{x}) + g(\mathbf{x})\nabla f(\mathbf{x})$
- **Chain rule:**
 - If $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and $g: \mathbb{R}^m \rightarrow \mathbb{R}^n$, then
 $\nabla(f \circ g)(\mathbf{x}) = \mathbf{J}(g(\mathbf{x}))^T (\nabla f(\mathbf{y}))$ where $\mathbf{y} = g(\mathbf{x})$
 - If f is as above and $h: \mathbb{R} \rightarrow \mathbb{R}$, then
 $\nabla(h \circ f)(\mathbf{x}) = h'(f(\mathbf{x}))\nabla f(\mathbf{x})$

IMPORTANT!

Convexity



- A function is **convex** if any line segment between two points of the function lie above or on the graph
- For univariate f , if $f''(x) \geq 0$ for all x
- For multivariate f , if its Hessian is **positive semidefinite**
 - I.e. $\mathbf{z}^T \mathbf{H} \mathbf{z} \geq 0$ for any \mathbf{z}
- Convex function's local minimum is its global minimum

Preserving the convexity

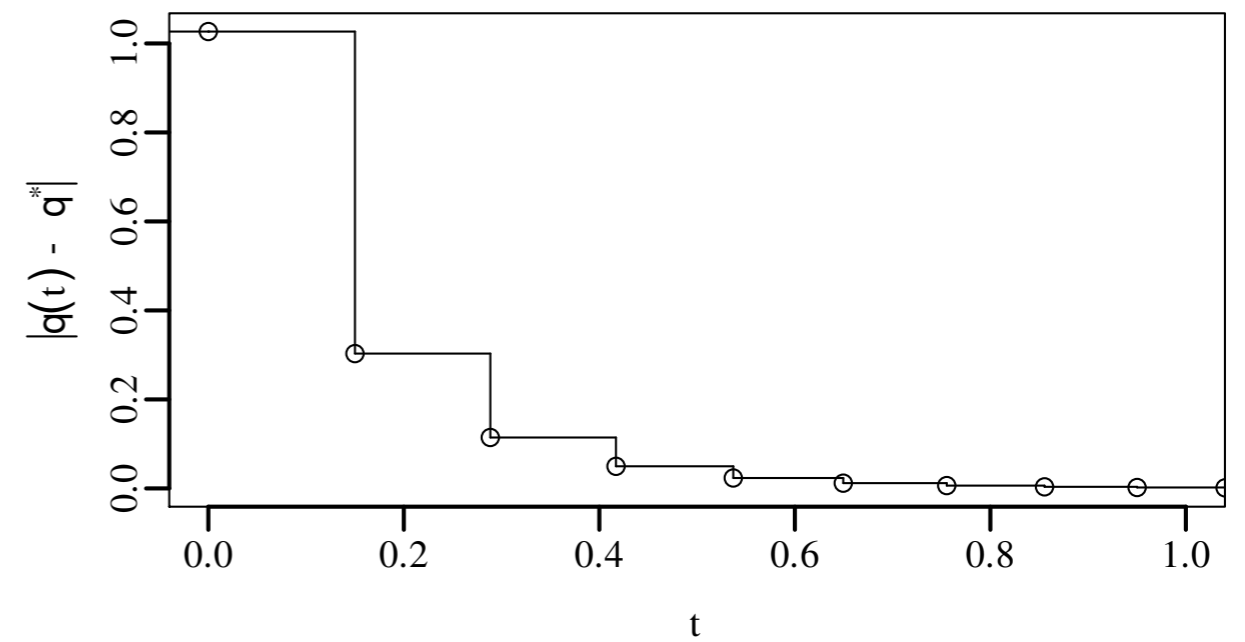
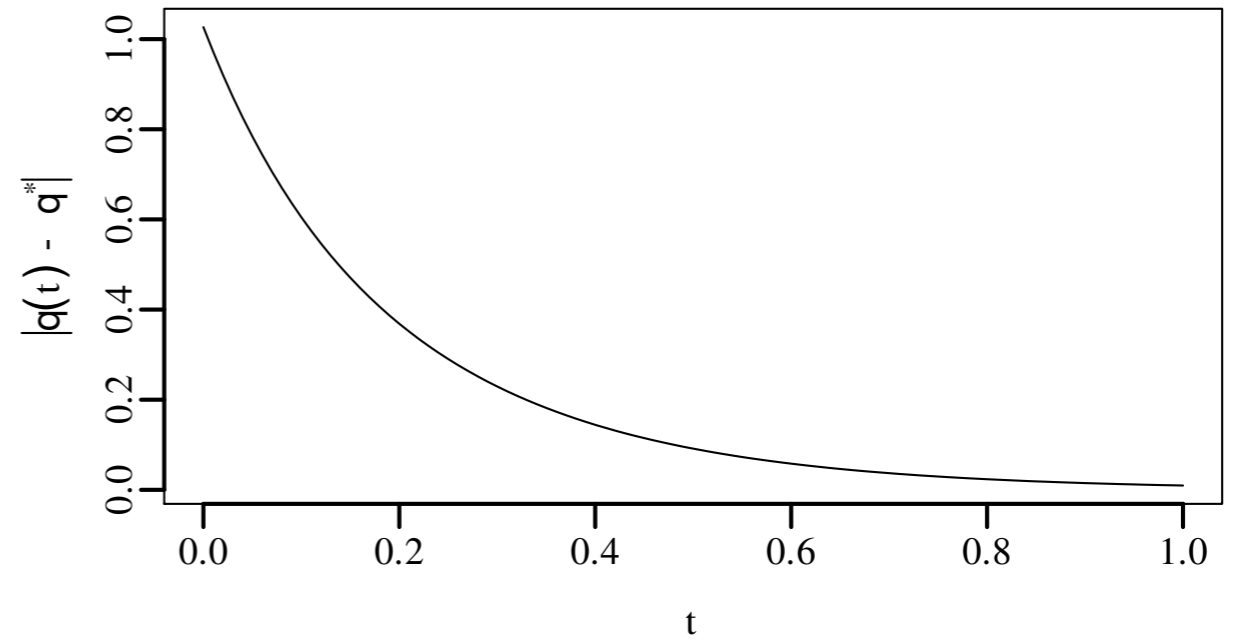
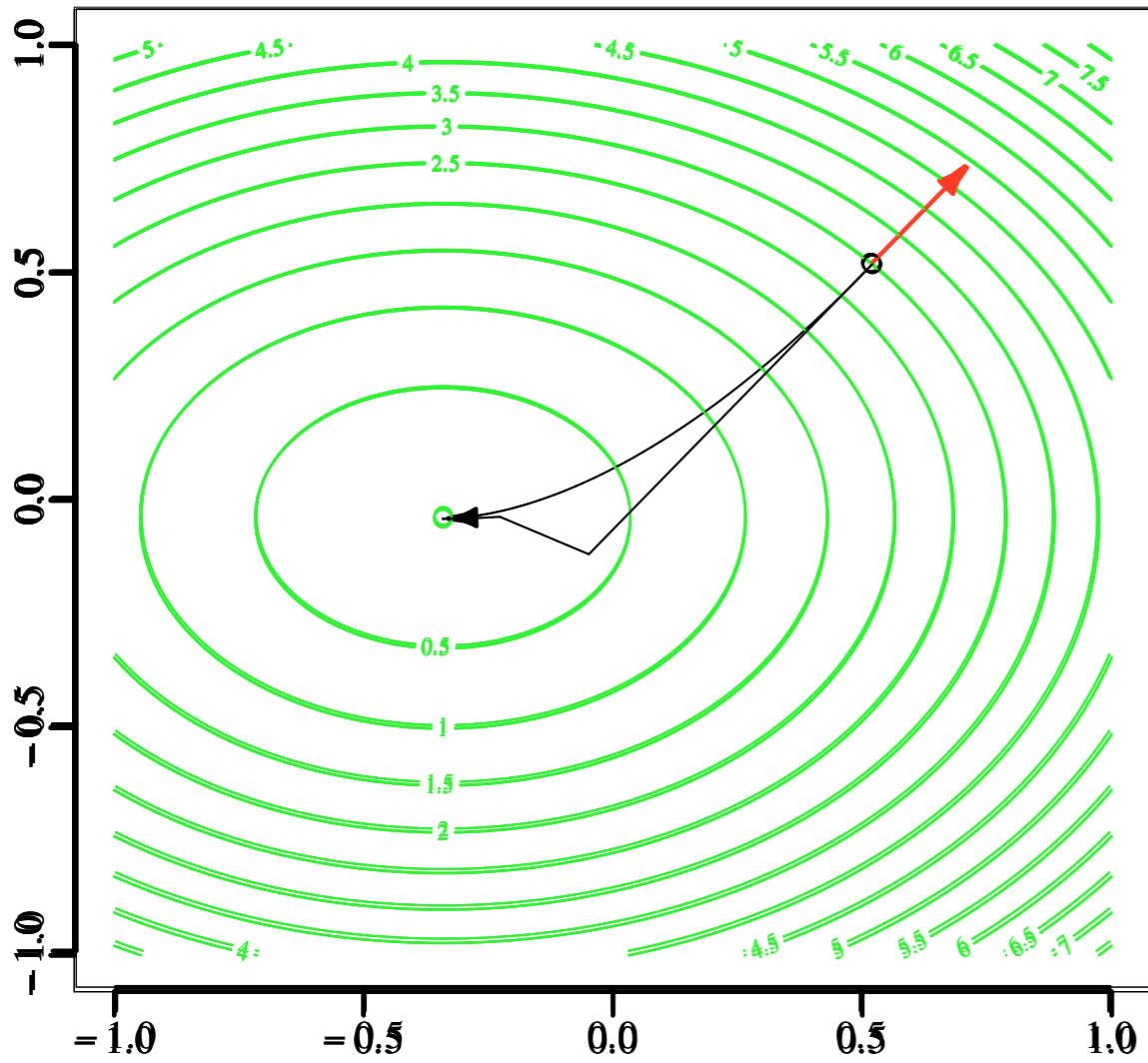
- If f is convex and $\lambda > 0$, then λf is convex
- If f and g are convex, the $f + g$ is convex
- If f is convex and g is **affine** (i.e. $g(\mathbf{x}) = \mathbf{Ax} + \mathbf{b}$), then $f \circ g$ is convex (N.B. $(f \circ g)(\mathbf{x}) = f(\mathbf{Ax} + \mathbf{b})$)
- Let $f(\mathbf{x}) = (h \circ g)(\mathbf{x})$ with $g: \mathbb{R}^n \rightarrow \mathbb{R}$ and $h: \mathbb{R} \rightarrow \mathbb{R}$; f is convex if
 - g is convex and h is nondecreasing and convex
 - g is concave and h is non-increasing and convex

Gradient descent

Idea

- If f is convex, we should find its minimum by following its negative gradient
- But the gradient at \mathbf{x} points to minimum only at \mathbf{x}
- Hence, we need to descent slowly down the gradient

Example



Gradient descent

- Start from random point \mathbf{x}^0
- At step n , update $\mathbf{x}^n \leftarrow \mathbf{x}^{n-1} - \gamma \nabla f(\mathbf{x}^{n-1})$
 - γ is some small **step size**
 - Often, γ depends on the iteration
$$\mathbf{x}^n \leftarrow \mathbf{x}^{n-1} - \gamma_n \nabla f(\mathbf{x}^{n-1})$$
- With suitable f and step size, will converge to local minimum

Example: least squares

- Given $\mathbf{A} \in \mathbb{R}^{n \times m}$ and $\mathbf{b} \in \mathbb{R}^n$, find $\mathbf{x} \in \mathbb{R}^m$
s.t. $\|\mathbf{Ax} - \mathbf{b}\|^2/2$ is minimized
 - Can be solved using SVD...
- Calculate the gradient of $f_{\mathbf{A},\mathbf{b}}(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|^2/2$
- Employ the gradient descent approach
 - In this case, the step size can be calculated analytically

Example: the gradient

Let's write open:

$$\begin{aligned}\frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 &= \frac{1}{2} \sum_{i=1}^n ((\mathbf{Ax})_i - b_i)^2 \\ &= \frac{1}{2} \sum_{i=1}^n \left(\sum_{j=1}^m a_{ij} x_j - b_i \right)^2 \\ &= \frac{1}{2} \sum_{i=1}^n \left(\left(\sum_{j=1}^m a_{ij} x_j \right)^2 - 2b_i \sum_{j=1}^m a_{ij} x_j + b_i^2 \right) \\ &= \frac{1}{2} \sum_{i=1}^n \left(\sum_{j=1}^m a_{ij} x_j \right)^2 - \sum_{i=1}^n b_i \sum_{j=1}^m a_{ij} x_j + \frac{1}{2} \sum_{i=1}^n b_i^2\end{aligned}$$

Example: the gradient

The partial derivative w.r.t. x_j :

$$\begin{aligned}
 \frac{\partial}{\partial x_j} \left(\frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 \right) &= \frac{\partial}{\partial x_j} \left(\frac{1}{2} \sum_{i=1}^n \left(\sum_{k=1}^m a_{ik} x_k \right)^2 - \sum_{i=1}^n b_i \sum_{k=1}^m a_{ik} x_k + \frac{1}{2} \sum_{i=1}^n b_i^2 \right) \\
 \text{Linearity} &= \frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial x_j} \left(\sum_{k=1}^m a_{ik} x_k \right)^2 - \sum_{i=1}^n b_i \frac{\partial}{\partial x_j} \sum_{k=1}^m a_{ik} x_k + \frac{\partial}{\partial x_j} \frac{1}{2} \sum_{i=1}^n b_i^2 \\
 &= \frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial x_j} \left(\sum_{k=1}^m a_{ik} x_k \right)^2 - \sum_{i=1}^n b_i a_{ij} \quad \begin{matrix} = 0 \text{ if } k \neq j \\ = 0 \end{matrix} \\
 \text{Chain rule} &= \sum_{i=1}^n a_{ij} \sum_{k=1}^m a_{ik} x_k - \sum_{i=1}^n b_i a_{ij} \\
 &= \sum_{i=1}^n a_{ij} \left(\sum_{k=1}^m a_{ik} x_k - b_i \right)
 \end{aligned}$$

Example: the gradient

Collecting terms:

$$\begin{aligned}\frac{\partial}{\partial x_j} \left(\frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 \right) &= \sum_{i=1}^n a_{ij} \left(\sum_{k=1}^m a_{ik} x_k - b_i \right) \\ &= \sum_{i=1}^n a_{ij} \left((\mathbf{Ax})_i - b_i \right) \\ &= \left(\mathbf{A}^T (\mathbf{Ax} - \mathbf{b}) \right)_j\end{aligned}$$

Matrix product

Another matrix product

Hence we have:

$$\nabla \left(\frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 \right) = \mathbf{A}^T (\mathbf{Ax} - \mathbf{b})$$

Example: the gradient

The other way: Use the chain rule

$$\begin{aligned}\nabla\left(\frac{1}{2}\|\mathbf{Ax} - \mathbf{b}\|^2\right) &= \mathbf{J}(\mathbf{Ax} - \mathbf{b})^T \left(\nabla\left(\frac{1}{2}\|\mathbf{y}\|^2\right)\right) \quad \mathbf{y} = \mathbf{Ax} - \mathbf{b} \\ &= \mathbf{A}^T (\mathbf{Ax} - \mathbf{b})\end{aligned}$$

Gradient descent & matrices

- How about “Given \mathbf{A} , find small \mathbf{B} and \mathbf{C} s.t. $\|\mathbf{A} - \mathbf{BC}\|_F$ is minimized”?
 - Not convex for \mathbf{B} and \mathbf{C} jointly
- Fix some \mathbf{B} and solve for \mathbf{C}
 - $\mathbf{C} = \operatorname{argmin}_{\mathbf{X}} \|\mathbf{A} - \mathbf{BX}\|_F$
- Use the found \mathbf{C} and solve for \mathbf{B} , and repeat until convergence

How to solve for C ?

- $C = \operatorname{argmin}_X ||\mathbf{A} - \mathbf{B}\mathbf{X}||_F$ still needs some work
- Write the norm as sum of column-wise errors

$$||\mathbf{A} - \mathbf{B}\mathbf{X}||_F = \sum ||\mathbf{a}_j - \mathbf{B}\mathbf{x}_j||_2$$

- Now the problem is a series of standard least-squares problems
- Each can be solved independently

How to select the step size?

- Recall: $\mathbf{x}^n \leftarrow \mathbf{x}^{n-1} - \gamma_n \nabla f(\mathbf{x}^{n-1})$
- Selecting correct γ_n for each n is crucial
 - Methods for optimal step size are often slow (e.g. line search)
 - Wrong step size can lead to non-convergence

Stochastic gradient descent

Basic idea

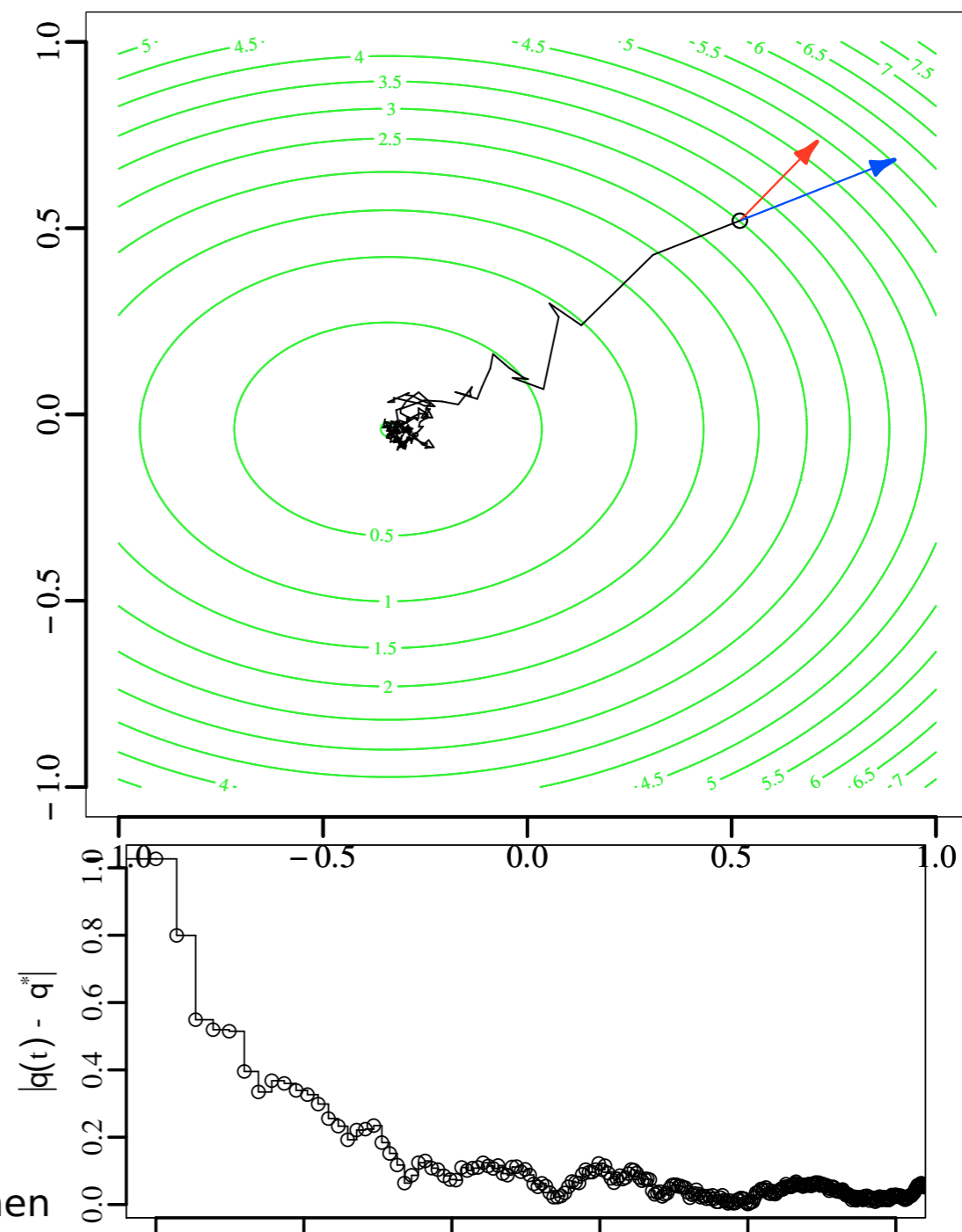
- With gradient descent, we need to calculate the gradient for $\mathbf{c} \mapsto \|\mathbf{a} - \mathbf{B}\mathbf{c}\|$ many times for different \mathbf{a} in each iteration
- Instead we can fix one element a_{ij} and update the i th row of \mathbf{B} and j th column of \mathbf{C} accordingly
- When we choose a_{ij} randomly, this is **stochastic gradient descent** (SGD)

Local gradient

- With fixed a_{ij} , $\|a_{ij} - (\mathbf{BC})_{ij}\| = a_{ij} - \sum b_{ik}c_{kj}$
 - Local gradient for b_{ik} is $-2c_{kj}(a_{ij} - (\mathbf{BC})_{ij})$
 - Similarly for c_{kj}
- This allows us to update the factors by only computing one gradient
 - Gradient needs to be sufficiently scaled

SGD process

- Initialize with random \mathbf{B} and \mathbf{C}
- **repeat**
 - Pick a random element (i, j)
 - Update a row of \mathbf{B} and a column of \mathbf{C} using the local gradients w.r.t. a_{ij}



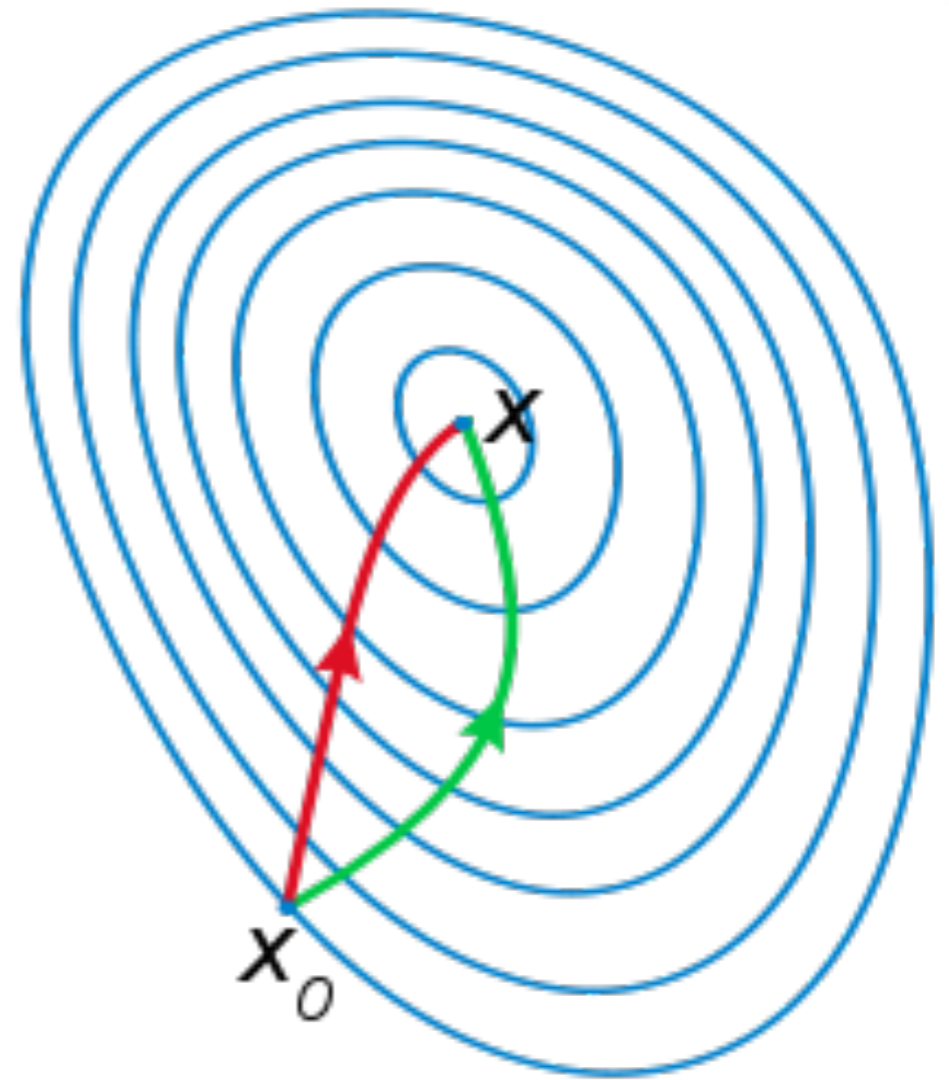
SGD pros and cons

- Each iteration is faster to compute
 - But can increase the error
- Does not need to know all elements of the input data
 - Scalability
 - Partially observed matrices (e.g. collaborative filtering)
- The step size still needs to be chosen carefully

Newton's method

Basic idea

- Iterative update rule:
$$\mathbf{x}_{n+1} \leftarrow \mathbf{x}_n - [\mathbf{H}(f(\mathbf{x}_n))]^{-1} \nabla f(\mathbf{x}_n)$$
- Assuming Hessian exists and is invertible...
- Takes curvature information into account



Pros and cons

- Much faster convergence
 - But Hessian is slow to compute and takes lots of memory
- Quasi-Newton methods (e.g. L-BFGS) compute the Hessian indirectly
- Often still needs some step size other than 1

Alternating least squares

Basic idea

- Given \mathbf{A} and \mathbf{B} , we can find \mathbf{C} that minimizes $\|\mathbf{A} - \mathbf{BC}\|_F$
- In gradient descent, we move slightly towards \mathbf{C}
- In **alternating least squares** (ALS), we replace \mathbf{C} with the new one

Basic ALS algorithm

- Given \mathbf{A} , sample a random \mathbf{B}
- **repeat** until convergence
 - $\mathbf{C} \leftarrow \operatorname{argmin}_{\mathbf{X}} \|\mathbf{A} - \mathbf{BX}\|_F$
 - $\mathbf{B} \leftarrow \operatorname{argmin}_{\mathbf{X}} \|\mathbf{A} - \mathbf{XC}\|_F$

ALS pros and cons

- Can have faster convergence than gradient descent (or SGD)
- The update is slower to compute than in SGD
 - About as fast as in gradient descent
- Requires fully-observed matrices

Adding constraints

The problem setting

- So far, we have done **unconstrained optimization**
- What if we have constraints on the optimal solution?
 - E.g. all matrices must be nonnegative
- In general, the above approaches won't admit these constraints

General case

- Minimize $f(\mathbf{x})$
- Subject to
$$g_i(\mathbf{x}) \leq 0, i = 1, \dots, m$$
$$h_j(\mathbf{x}) = 0, j = 1, \dots, k$$
- Assuming certain regularity conditions, there exists constraints μ_i ($i=1, \dots, m$) and λ_j ($j=1, \dots, k$) that satisfy **Karush–Kuhn–Tucker** (KKT) conditions

KKT conditions

- Let \mathbf{x}^* be the optimal solution
- **Stationarity:**
 - $-\nabla f(\mathbf{x}^*) = \sum_i \mu_i \nabla g_i(\mathbf{x}^*) + \sum_j \lambda_j \nabla h_j(\mathbf{x}^*)$
- **Primal feasibility:**
 - $g_i(\mathbf{x}^*) \leq 0$ for all $i = 1, \dots, m$
 - $h_j(\mathbf{x}^*) = 0$ for all $j = 1, \dots, k$
- **Dual feasibility:**
 - $\mu_i \geq 0$ for all $i = 1, \dots, m$
- **Complementary slackness:**
 - $\mu_i g_i(\mathbf{x}^*) = 0$ for all $i = 1, \dots, m$

When do KKT conditions hold

- KKT conditions hold under certain regularity conditions
 - E.g. g_i and h_j are affine
 - Or f is convex and exists \mathbf{x} s.t. $h(\mathbf{x}) = 0$ and $g_i(\mathbf{x}) < 0$
- Nonnegativity is an example of linear (hence, affine) constraint

What to do with the KKT conditions?

- μ and λ are new unknown variables
 - Must be optimized together with \mathbf{x}
- The conditions appear in the optimization
 - E.g. in the gradient
- The KKT conditions are rarely solved directly

Summary

- There are **many** methods for optimization
 - We only scratched the surface
- Methods are often based on gradients
 - Can lead into ugly equations
- Next week: applying these techniques for finding nonnegative factorizations... *Stay tuned!*